

```
from pprint import pprint
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfTransformer
```

```

from sklearn.naive_bayes import MultinomialNB
from sklearn import metrics
from sklearn.feature_extraction.text import TfidfVectorizer
from nltk.tokenize import sent_tokenize, word_tokenize
import warnings
from gensim.models import Word2Vec
from gensim.test.utils import common_texts
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.datasets import fetch_20newsgroups
from sklearn.model_selection import train_test_split
from time import time
from sklearn.model_selection import GridSearchCV
from sklearn.tree import DecisionTreeClassifier
from sklearn.svm import SVC
from sklearn.naive_bayes import MultinomialNB
from multiprocessing import cpu_count
from sklearn import feature_selection
from keras.preprocessing.text import Tokenizer, text_to_word_sequence
from keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.utils import to_categorical
from keras.models import Model
from keras.layers import Dense, Input, Flatten, Reshape, concatenate, Dropout
from keras.layers import Conv1D, Conv2D, MaxPooling1D, MaxPooling2D, Embedding
from keras import optimizers
from keras.callbacks import EarlyStopping
import tensorflow as tf

# context
newsgroups_train = fetch_20newsgroups(subset='train', random_state=42, shuffle=True)
print(list(newsgroups_train.target_names))

('alt.atheism',
 'comp.graphics',
 'comp.os.msc-windows.misc',
 'comp.sys.ibm.pc.hardware',
 'comp.sys.mac.hardware',
 'comp.windows.x',
 'misc.forsale',
 'rec.autos',
 'rec.motorcycles',
 'rec.sport.baseball',
 'rec.sport.hockey',
 'sci.crypt',
 'sci.electronics',
 'sci.med',
 'sci.space',
 'soc.religion.christian',
 'talk.politics.guns',
 'talk.politics.mideast',
 'talk.politics.misc',
 'talk.religion.misc')

# size
len(newsgroups_train.data)
len(newsgroups_train filenames)
print(f"Number of files: {len(newsgroups_train.filenames)}")

print("\n".join(newsgroups_train.data[0].split("\n")[:3]))
print(newsgroups_train.target_names[newsgroups_train.target[0]])

print(newsgroups_train.target[10])
for t in newsgroups_train.target[:10]:
    print(newsgroups_train.target_names[t])
    print("====/====/====/====")

for name in newsgroups_train.target_names:
    news = fetch_20newsgroups(subset='train', categories=[name], data
    print(f"Category {name} has {len(news)} records")

Number of files: 11314
From: lerkst@wan.umd.edu (where's my thing)
Subject: WHAT car is this?
http://posting-host: rac).wan.umd.edu
rec.autos
7 7 4 1 14 16 13 3 2 4]
rec.autos
comp.sys.mac.hardware
comp.sys.mac.hardware
comp.graphics
sci.space
talk.politics.guns
sci.med
comp.sys.ibm.pc.hardware
comp.os.ms-windows.misc
comp.sys.mac.hardware
====/====/====/====
Category alt.atheism has 480 records
Category comp.graphics has 384 records
Category comp.os.ms-windows.misc has 591 records
Category comp.sys.ibm.pc.hardware has 590 records
Category comp.sys.mac.hardware has 578 records
Category comp.windows.x has 593 records
Category misc.forsale has 585 records
Category rec.autos has 594 records
Category rec.motorcycles has 598 records
Category rec.sport.baseball has 597 records
Category rec.sport.hockey has 600 records
Category sci.crypt has 595 records
Category sci.electronics has 591 records
Category sci.med has 594 records
Category sci.space has 593 records
Category soc.religion.christian has 599 records
Category talk.politics.guns has 546 records
Category talk.politics.mideast has 564 records
Category talk.politics.misc has 465 records
Category talk.religion.misc has 377 records

difficulties and objectives? In progress

1. Assign a fixed integer id to each word occurring in any document of the training set (for instance by building a dictionary from words to integer indices).

Below I am trying to modify the data and to use bag of words. I'll start with lowering the case of all text and ignoring/removing frequent words.

# print("\n".join(newsgroups_train.data[0].split("\n")[:3]))
common_words = ['the', 'at', 'there', 'some', 'my', 'of', 'be', 'use', 'her', 'than', 'and', 'this', 'an', 'would', 'first',
                'never', 'to', 'from', 'which', 'like', 'been', 'in', 'for', 'she', 'him', 'all', 'is', 'one', 'do', 'into',
                'time', 'oil', 'that', 'by', 'their', 'has', 'its', 'it', 'word', 'if', 'look', 'now', 'be', 'but', 'will',
                'up', 'more', 'long', 'for', 'what', 'other', 'write', 'down', 'on', 'all', 'about', 'go', 'day', 'are', 'we',
                'a', 'we', 'many', 'number', 'get', 'with', 'when', 'then', 'no', 'come', 'his', 'your', 'them', 'way', 'm

data for groups = {}
newsgroups_train["data"]

```

```

new_data['word'] = filtered_data
filtered_data = [re.sub('\{', '', new) for new in new_data]
filtered_data = [re.sub('\}', '', new) for new in filtered_data]
filtered_data = [re.sub('\^', '', new) for new in filtered_data]
filtered_data = [re.sub('\$', '', new) for new in filtered_data]
filtered_data = [re.sub('\[', '', new) for new in filtered_data]
filtered_data = [re.sub('\]', '', new) for new in filtered_data]
for word in common_words:
    filtered_data = [new.replace(' ' + word + ' ', ' ') for new in filtered_data]
filtered_data = [new.lower() for new in filtered_data]
for d in filtered_data:
    d = ' '.join([w for w in d.split() if w != ''])
# filtered_data = [new.strip() for new in filtered_data]
# didnt manage to make it work in time
return filtered_data

# filtered_data = [re.sub(r'["\w] ', ' ', new) for new in new_data]
# for word in common_words:
#     filtered_data = [new.replace(word, ' ') for new in filtered_data]
# filtered_data = [new.lower() for new in filtered_data]
# filtered_data = [new.lower() for new in filtered_data]
# filtered_data should be valid
filtered_data = filter_text_data(new_data)

for name, data in zip(newsgroups_train.target_names, newsgroups_train.data):
    info_about_group = {}
    info_about_group['name'] = name
    info_about_group['word_dict'] = {}
    info_about_group['word_occurrence'] = {}
    info_about_group['word_occurrence_indexes_and_occurrences'] = {}
    data = re.sub('["\w] ', ' ', data)
    data = data.lower()
    integer = 0

reversed_word_occurrence_dict = {}
for word in data.split(' '):
    if word and len(word)>1:
        if word not in info_about_group['word_dict'] and word not in common_words:
            info_about_group['word_dict'][word] = integer
            integer += 1
for word in info_about_group['word_dict']:
    info_about_group['word_occurrence'][word] = 0

for word in data.split(' '):
    if word and len(word)>1:
        if word in info_about_group['word_occurrence']:
            info_about_group['word_occurrence'][word] += 1
data_for_groups.append(info_about_group)

for integer, occurrence in zip(info_about_group['word_dict'].values(), info_about_group['word_occurrence'].values()):
    info_about_group['word_occurrence_indexes_and_occurrences'][integer] = occurrence
print(data_for_groups[4]['word_occurrence_indexes_and_occurrences'])

(0, 1), 1, 2, 1, 3, 1, 4, 3, 5, 2, 6, 3, 7, 1, 8, 1, 9, 1, 10, 2, 11, 2, 12, 1, 13, 1, 14, 1, 15, 1, 16, 1, 1, 17, 1, 18, 1, 19, 1, 20, 1, 21, 1, 22, 1, 23, 2, 24, 1, 25, 1, 26, 2, 27, 2, 28, 2, 29, 1, 30, 2, 31, 1, 32, 1, 33, 1, 34, 1, 35, 1, 36, 1, 37, 1, 38, 1, 39, 1, 40, 1, 41, 1, 42, 1, 43, 1, 44, 1, 45, 2, 46, 1, 47, 1, 48, 2, 49, 1, 50, 1, 51, 2, 52, 1, 53, 1, 54, 1, 55, 1, 56, 1, 57, 1, 58, 1, 59, 2, 60, 1, 61, 1, 62, 1, 63, 1, 64, 1, 65, 1, 66, 1, 67, 1, 68, 1, 69, 1, 70, 1, 71, 2, 72, 1, 73, 1, 74, 1, 75, 1, 76, 1, 77, 1, 78, 1, 79, 1, 80, 1, 81, 1, 82, 1, 83, 1, 84, 1, 85, 1, 86, 1, 87, 1, 88, 1, 89, 1, 90, 1, 91, 1, 92, 1, 93, 1, 94, 1, 95, 1, 96, 1, 97, 1, 98, 1, 99, 1, 100, 1)

# I guess we are counting words now with some pre-built functions
count_vect = CountVecTransformer()
X_train_counts = count_vect.fit_transform(filtered_data)
print(X_train_counts.shape)
print(count_vect.vocabulary_.get('ulgorith'))

(11314, 130107)
27366

# as my friend suggested to try trid transform as well, because words frequency can be better than just word
tridf_transformer = TridfTransformer()
X_train_tridf = tridf_transformer.fit_transform(X_train_counts)
# This is an example with formatted data
tridf_vectorizer = TridfVectorizer()
vectors = tridf_vectorizer.fit_transform(filtered_data)

# So, if we use already defined functions to vectorize, fit, predict, here it is
newsgroups_test = fetch_20newsgroups(subset='test', random_state=42, shuffle=True)
```

```
clf = MultinomialNB(alpha=.01)
clf.fit(newgroups, newgroups_train.target)
pred = clf.predict(newgroups_test)
metrics.f1_score(newgroups_test.target, pred, average='macro')

0.8299999321458893

# When the torture is done, let's stick to predefined functions
# this is an example with original data
vectorizer = TfidfVectorizer()
newgroups = vectorizer.fit_transform(newgroups_train.data)
# we also use predefined functions for vectorizer, fit and predict, here it is
newgroups_test = fetch_20newsgroups(subset='test', random_state=42, shuffle=True)
newgroups_test = vectorizer.transform(newgroups_test.data)
clf = MultinomialNB(alpha=.01)
```

```
pred = clf.predict(vectors_test)
metrics: f1_score(newsgroups_test.target, pred, average='macro')
Overall, Not worth the fuss
0.6290659644474043
```

```
EMBEDDING_DIM = 100
MAX_SEQUENCE_LENGTH = 1000
MAX_NUM_WORDS = 20000

def text_cnn(embedding_layer):
    sequence_input = Input(shape=(MAX_LEN,), dtype='int32')
    embedded_sequences = embedding_layer(sequence_input)

    # Yoon Kim model (https://arxiv.org/abs/1408.5882)
    embedded_sequences = Reshape((MAX_LEN, EMBEDDING_DIM, 1))(embedded_sequences)
    x = Conv2D(100, (5, EMBEDDING_DIM), activation='relu')(embedded_sequences)
    x = MaxPooling2D((MAX_LEN - 5 + 1, 1))(x)

    y = Conv2D(100, (4, EMBEDDING_DIM), activation='relu')(embedded_sequences)
```

```
z = Conv2D(100, (3, 3), embedding_dim, activation='relu')(embedded_sequences)
z = MaxPooling2D(MAX_DIM - 3 + 1, 1)(z)

alpha = concatenate([w,y,z])
alpha = Flatten()(alpha)
alpha = Dropout(0.5)(alpha)
preds = DenseLayer(newsgroups_train.target_names, activation='softmax')(alpha)
model = Model(sequence_input, preds)
adadelta = tf.optimizers.Adadelta()

model.compile(loss='categorical_crossentropy',
               optimizer=adadelta,
               metrics=['acc'])

return model
```

```

> show.history(history):
plt.plot(history['acc'])
plt.plot(history['val_acc'])
plt.title('model accuracy')
plt.ylabel('accuracy')
plt.xlabel('epoch')
plt.legend(['train', 'test'], loc='upper left')
plt.show()

# summarize history for loss
plt.plot(history['loss'])
plt.plot(history['val_loss'])
plt.title('model loss')
plt.ylabel('loss')
plt.xlabel('epoch')
plt.legend(['train', 'test'], loc='upper left')

```

```
def show_performance_with_gscv(name, model, x_train, y_train, params):
    results = {}
    results['model_name'] = name
    gcv = GridSearchCV(model, params, cv=5, n_jobs=cpu_count()-1, return_train_score=True)
    # x_train_fs = fs.flit_transform(x_train, y_train)

    gscv.fit(x_train, y_train)

    results['params'] = gscv.best_params_
    results['train_time'] = np.mean(gscv.cv_results_['mean_fit_time'])
    results['val_time'] = np.mean(gscv.cv_results_['mean_score_time'])
    results['train_score'] = gscv.cv_results_['mean_train_score'][gscv.best_index_]
    results['val_score'] = gscv.cv_results_['mean_test_score'][gscv.best_index_]

    # it is get for train data set, could be taken as a val result
```

```
results['best_model'] = gsvc.best_estimator_

return results

def show_distribution(data):
    dict = {}
    for name, name in enumerate(newsgroups_train.target_names):
        dict.setdefault(name, np.sum(data==index))
    print(dict)
    print(dict.keys())
    print(dict.values())

    index = np.arange(len(newsgroups_train.target_names))
    plt.figure(figsize=(10,8))
    plt.bar(index, dict.values())
```

```
plt.title("category distribution")
plt.xlabel("data count")
plt.ylabel("data category")
plt.show()

def show_performance(model, x_train, y_train, x_val, y_val):
    results = {}
    results['model_name'] = model.__class__.__name__
    t0 = time()
    model.fit(x_train, y_train)
    results['train_time'] = time() - t0
    t1 = time()
    predicts = model.predict(x_train)
    results['val_time'] = time() - t1
    train_score = model.score(x_train, y_train)
```

```

results['train_score'] = train_score
results['val_score'] = val_score
print(results)

def show_words(data):
    count = {}
    for f in data:
        count.append(len(f.split()))
    plt.figure(figsize=(10,5))
    plt.hist(count, bins=20)
    plt.title("words distribution")
    plt.xlabel("words count")
    plt.ylabel("words weight")
    plt.show()

```

```
def show_chars(data):
    count = 1
    for f in data:
        count.append(len(f))
    plt.figure(figsize=(10,3))
    plt.hist(count, bins=20)
    plt.title("chars distribution")
    plt.xlabel("chars count")
    plt.ylabel("chars count")
    plt.show()

def autolabel(ax, rects, xpos='center'):
    xpos = xpos.lower() # normalize the case of the parameter
    y_pos = ['center', 'top', 'bottom']
    offset = {}
    offset = {'center': 0.5, 'right': 0.57, 'left': 0.43} # x_txt = x + w*off
```

```

for rect in rects:
    height = rect.get_height()
    ax.text(rect.get_x() + rect.get_width() * 0.5, height,
            '{}'.format(height, ha='right', va='bottom'))

def show_metrics(rs_lr, rs_svm, rs_nb):
    train_time = [rs_lr['train_time'], rs_svm['train_time'], rs_nb['train_time']]
    val_time = [rs_lr['val_time'], rs_svm['val_time'], rs_nb['val_time']]
    train_score = [rs_lr['train_score'], rs_svm['train_score'], rs_nb['train_score']]
    val_score = [rs_lr['val_score'], rs_svm['val_score'], rs_nb['val_score']]

    ind = np.arange(1, len(train_time)) # the x locations for the groups
    width = 0.35 # the width of the bars

    fig, ax = plt.subplots(1, 1, figsize=(16, 6))

```

```
rects1 = ax0.bar(ind - width/2, train_time, width,
                 color='skyblue', label='train time')
rects2 = ax0.bar(ind + width/2, val_time, width,
                 color='indianred', label='val time')

# Add some text for labels, title and custom x-axis tick labels, etc.
ax0.set_ylabel('time')
ax0.set_title('time by train and val')
ax0.set_xticks(ind)
ax0.set_xticklabels(("LogisticRegression", "SVC", "NaiveBayes"))
ax0.legend()
autolabel(ax0, rects1, 'left')
autolabel(ax0, rects2, 'right')
plt.show()
```

```
fig, axes = plt.subplots(2, 1, figsize=(10, 10))

rects3 = ax1.bar(ind - width/2, train_score, width,
                 color='blue', label='train score')
rects4 = ax1.bar(ind + width/2, val_score, width,
                 color='IndianRed', label='val score')
ax1.set_ylabel('Score')
ax1.set_title('Scores by train and val')
ax1.set_xticks(ind)
ax1.set_xticklabels(['LogisticRegression', 'SVC', 'NaiveBayes'])
ax1.legend()
autolabel(ax1, rects3, 'left')
autolabel(ax1, rects4, 'right')
plt.show()
```

```
new_groups_train = fetch_20newsgroups(subset='train', shuffle=True, download_if_missing=False)
new_groups_test = fetch_20newsgroups(subset='test', shuffle=True, download_if_missing=False)
new_group_train_filter_text_data = new_groups_train.data
new_group_train_filter_text_target = new_groups_train.target
new_group_test_filter_text_data = new_groups_test.data
new_group_test_filter_text_target = new_groups_test.target
x_train, y_train = filtered_train, new_group_train_target
x_test, y_test = filtered_train, new_group_test_target

print('Count for data in 20 news groups', len(x_sp_train), len(x_sp_val), len(x_sp_test))
print('Count for train and validation data in 20 news groups', len(x_sp_train) + len(x_sp_val), " and for test")
print('Count for train data in 20 news groups', len(x_sp_train))
print('Count for validation data in 20 news groups', len(x_sp_val))
print('Count for test data in 20 news groups', len(x_sp_test))
```

```
count for data in 20 news groups 9051 2863 13134
count for train data in 20 news groups 9051 2863 13134 and for test data 13134
count for train data in 20 news groups 9051
count for validation data in 20 news groups 2263
count for test data in 20 news groups 13134

print(x_sp_train[10])

from djf.cok.covcovary import uk_marvin_batby subject re moonbase race nntp posting host cc_sysk organization
startfilecount codelink uk_lines 22 in article 1469091nn14j nojo eng umd edu ysamgr link eng and ude writers
in article codelink 7q3 zo3 toronto ede h1r500004j nojo eng umd edu hspencer writers spollo dno hard bi
in article codelink 7q3 zo3 toronto ede h1r500004j nojo eng umd edu hspencer writers spollo dno hard bi
subject cuts costs factor several so much cost entirely venture assuming could talk u s government leasing
import cuts florida why ground launch pad it privately possible launch altitude this shuttle originally int
count for test data in 20 news groups 13134
```

```

      marvin base      djf uk ac cov cck and they shall those things sort rafia base fathers put
just night before at 8 o'clock

show_distribution(y_sp_train)

["'alt.theism': 392, 'comp.graphics': 456, 'comp.sys.windows.midea': 478, 'comp.sys.ibm.pc.hardware': 462, 'co
mp.sys.m.hardware': 'comp.windows.x': 473, 'misc.forsale': 480, 'rec.autos': 477, 'rec.motorcycles': 47
6, 'rec.sport.baseball': 478, 'rec.sport.hockey': 485, 'sci.crypt': 465, 'sci.electronics': 449, 'sci.education': 495, 'sci.gen
eral': 495, 'sci.space': 471, 'sci.religion.christian': 472, 'talk.politics.guns': 431, 'talk.politics.mideast': 452, 't
alk.politics.misc': 38, 'talk.religion.misc': 302]

['alt.theism', 'alt.theism', 'comp.graphics', 'comp.sys.windows.midea', 'comp.sys.ibm.pc.hardware', 'comp.sys.ma
.hardware', 'comp.windows.x', 'misc.forsale', 'rec.autos', 'rec.motorcycles', 'rec.sport.baseball', 'rec.spor
t.hockey', 'sci.electronics', 'sci.electronics', 'sci.space', 'sci.religion.christian', 'talk.politics.guns',
'talk.politics.mideast', 'talk.politics.misc', 'talk.religion.misc']

```

category	count
category 0	390
category 1	450
category 2	470
category 3	460
category 4	450
category 5	470
category 6	470
category 7	460
category 8	480
category 9	480
category 10	460
category 11	480
category 12	490
category 13	480
category 14	470
category 15	420
category 16	380
category 17	300
category 18	320
category 19	310

[illegible]

```

"cat.system.chairman": 88, "comp.graphics": 128, "comp.os.ms-windows.misc": 113, "comp.sys.ibm.pc.hardware": 128, "comp.sys.mac.hardware": 120, "comp.windows.x": 120, "misc.forsale": 105, "misc.hardware": 117, "rec.horticulture": 128, "rec.sport.baseball": 119, "rec.sport.hockey": 115, "sci.crypt": 130, "sci.electronics": 122, "sci.edu": 99, "sci.space": 122, "soc.religion.christian": 137, "talk.politics.guns": 115, "talk.politics.mideast": 112, "talk.politics.misc": 80, "talk.religion.misc": 75)
dict_keys(['cat.system.chairman', 'comp.graphics', 'comp.os.ms-windows.misc', 'comp.sys.ibm.pc.hardware', 'comp.sys.mac.hardware', 'comp.windows.x', 'misc.forsale', 'misc.hardware', 'rec.horticulture', 'rec.sport.baseball', 'rec.sport.hockey', 'sci.crypt', 'sci.electronics', 'sci.edu', 'sci.space', 'soc.religion.christian', 'talk.politics.guns', 'talk.politics.mideast', 'talk.politics.misc', 'talk.religion.misc'])
dict_values([88, 128, 113, 128, 120, 105, 117, 128, 115, 119, 122, 99, 122, 115, 112, 80, 75])

```

[illegible][illegible]

```

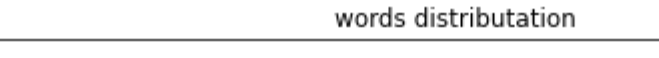
dic_keys["ai.theism", "comp.graphics", "comp.os.ms-windows.misc", "comp.sys.ibm.pc.hardware", "comp.sys.ni
comp.windows.x", "misc.forsale", "rec.autos", "rec.motorcycles", "rec.sport.baseball", "rec.spor
theology", "sci.crypt", "sci.electronics", "sci.med", "sci.space", "soc.religion.christian", "talk.politics.gun
", "talk.politics.mideast", "talk.politics.misc", "talk.religion.misc"]]
dic_values[[319, 389, 394, 392, 385, 395, 390, 396, 397, 395, 396, 393, 394, 398, 364, 376, 310, 25
]]

```

Topic	Number of Documents
all religion	145
comp graphics	140
cor ms windows misc	135
sys den pt hardware	130
sys den pt hardware	125
comp windows x	120
misc forscale	115
rec audio	110
rec motorcy cycle	105
rec sport basketball	100
rec sport hockey	95
sci crypt	90
sci electronics	85
sci med	80
sci space	75
soc religion christian	70
talk politics guns	65
talk politics midwest	60
talk politics misc	55
talk religion misc	50

```
data.count()
```

```
show_words(x_sp_train)
```



The chart displays the frequency of words in the training set. The most frequent word is 'the', with a count of approximately 8500. Other words like 'and', 'of', 'a', 'in', and 'is' also have high frequencies, around 7000 each. The chart shows a long tail of words with decreasing frequencies.

words count

count

show_charts(x_cp_train)

chars distribution

COUNTRY	OTHER COUNTRIES
Canada	8000
China	1000
France	500
Germany	500
India	500
Italy	500
Japan	500
Korea	500
Mexico	500
Russia	500
South Africa	500
Sweden	500
Switzerland	500
Taiwan	500
Thailand	500
United Kingdom	500
USA	8500

```
vectorizer = CountVectorizer(max_df=0.97, min_df=3,
                             max_features=None,
                             stop_words='english')
vec_x_train = vectorizer.fit_transform(x_sp_train)
print(vec_x_train)
```

```
(0, 3690) 5
(0, 11963) 1
(0, 31540) 1
(0, 27447) 2
(0, 27448) 2
```

(0, 29875)	1
(0, 4835)	3
(0, 9808)	2
(0, 843)	1
(0, 22761)	1
(0, 3992)	1
(0, 8594)	1
(0, 18156)	1
(0, 15031)	2
(0, 21245)	1
(0, 2317)	1
(0, 1364)	1
(0, 1229)	1
(0, 21345)	2
(0, 27183)	1
(0, 35010)	2

```

(0, 29670) 1
(0, 22295) 1
(0, 21824) 1
      1
(9050, 4638) 1
(9050, 10294) 1
(9050, 19577) 1
(9050, 30051) 1
(9050, 10748) 1
(9050, 19840) 1
(9050, 6631) 1
(9050, 27989) 1
(9050, 7654) 1
(9050, 21566) 1
(9050, 55071) 1

```

```

(9050, 3705) 1
(9050, 22290) 1
(9050, 1942) 3
(9050, 32680) 2
(9050, 8257) 1
(9050, 23393) 1
(9050, 1058) 1
(9050, 28660) 1
(9050, 30373) 1
(9050, 19568) 1
(9050, 18964) 1
(9050, 2639) 1
(9050, 29399) 1

```

```
stop_words='english')

tfidf_x_sp_train = tfidf_vectorizer.fit_transform(x_sp_train)
tfidf_x_sp_val = tfidf_vectorizer.transform(x_sp_val)

tfidf_x_train = tfidf_vectorizer.transform(x_train)
tfidf_x_test = tfidf_vectorizer.transform(x_test)

print(tfidf_x_sp_train.shape)
print(tfidf_x_sp_val.shape)
print(tfidf_x_train.shape)
print(tfidf_x_test.shape)

print(tfidf_x_sp_train)
```

```

(2263, 492504)
(11314, 492504)
(11334, 492504)
(0, 9116)
(0, 34277)
(0, 25839)
(0, 13630)
(0, 29653)
(0, 47241)
(0, 48099)
(0, 42799)
(0, 19402)
(0, 40026)
(0, 17118)
(0, 272)
(0, 056491052833452474
(0, 04662122634448474
(0, 02907127436323408
(0, 041713413543167023
(0, 03179398083128389)
(0, 0417726629923154
(0, 03232820196590292
(0, 02561667902683956
(0, 0361332703224655
(0, 04463461887716532
(0, 04231944622189823
(0, 03002167402050087

```

```
(0, 4910), 0.04515399280911522
(0, 711), 0.028188317553130602
(0, 7225), 0.04740972401374874
(0, 48744), 0.06570650973827732
(0, 27666), 0.06370650973827732
(0, 23140), 0.045103001343693725
(0, 32647), 0.04640943771810875
(0, 10028), 0.04008333487239442
(0, 38364), 0.039033810037373964
(0, 38113), 0.04941528920460408
(0, 44356), 0.02648045584629383
(0, 10580), 0.048621226384688774
:
(9050, 23458), 0.09724483693793684
(9050, 34208), 0.09556894303391444
```

```
(9050, 39479) 0.06781027097933104
(9050, 85357) 0.11333879680403721
(9050, 19331) 0.06618916014778352
(9050, 9693) 0.1749338047661834
(9050, 13350) 0.0696332407918378
(9050, 8760) 0.0745387959082176
(9050, 12840) 0.0957956696872671
(9050, 45312) 0.11115512707726367
(9050, 11443) 0.077420516182295792
(9050, 955) 0.06618916014778352
(9050, 25661) 0.08587668813878441
(9050, 36577) 0.10481938305900954
(9050, 15292) 0.08641994913598402
(9050, 21419) 0.113946485471834
(9050, 26764) 0.04278601325232935
```

```
(9050, 35003) 0.06123371207051186
(9050, 48712) 0.04089357747474934
(9050, 14576) 0.1169052846737264
(9050, 48069) 0.03198275590425969
(9050, 17617) 0.054642461185993229
(9050, 33136) 0.01991836998224083

# make word2vec model
model = Word2Vec(sentences=common_texts, vector_size=100, window=5, min_count=1, workers=4)
model.save("word2vec.model")

lr = LogisticRegression(C=0.1, penalty='l2')
show_performance(lr, tfidf_x_sp_train, y_sp_train, tfidf_x_sp_val, y_sp_val)
```

```

'model_name': 'LogisticRegression', 'train_time': 13.85401181030273, 'val_time': 0.01900005340576172, 'train_score': 0.9791132346663574, 'val_score': 0.8952487461985934)

svc = SVC(kernel='linear', C=0.5, gamma=0, random_state=0)
show_performance(svc, tfidf_x_sp_train, y_sp_train, tfidf_x_sp_val, y_sp_val)

'model_name': 'SVC', 'train_time': 6.7835351234741, 'val_time': 55.86075687408447, 'train_score': 0.98243288
4347133, 'val_score': 0.9032258064516219)

gnb = MultinomialNB(alpha=0.5)
show_performance(gnb, tfidf_x_sp_train, y_sp_train, tfidf_x_sp_val, y_sp_val)

'model_name': 'MultinomialNB', 'train_time': 0.04699373245329258, 'val_time': 0.02100467681847656, 'train_score':
0.9712738923878184, 'val_score': 0.89136632788334)

```

```
params = {'C': [0.01, 1, 3]}
rars_lm = show_performance_with_gscv('LogisticRegression', LogisticRegression(penalty='l2'), tfidf_x_train, y_train,
                                     list_train_ids)
{'model_name': 'LogisticRegression', 'params': {'C': 3}, 'train_time': 45.50942958725824, 'val_time': 0.0424789322161245, 'train_score': 0.995359719569596, 'val_score': 0.9040129370025266, 'best_model': 'LogisticRegression'}
/Users/roberta/AppData/Local/Programs/Python/Python39\lib/site-packages/sklearn/linear_model/_logistic.py:814:
ConvergenceWarning: LBFGS failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
```

```

n_iter_1 = check_optimize_result(
    params = {'C': [0.1, 1, 3], 'gamma': [0.5, 0.9]}
    rb_svc = show_performance_with_gsvc('SVC', SVC(kernel='linear'), tfidf_x_train, y_train, params)
    print(rb_svc)

    {'model_name': 'SVC', 'params': {'C': 3, 'gamma': 0.5}, 'train_time': 111.33690526705352, 'val_time': 34.8662720
    7264648, 'train_score': 0.9992487223052665, 'val_score': 0.906753314195242, 'best_model': SVC(C=3, gamma=0.5,
    kernel='linear')}

    params = {'alpha': [0.0001, 0.01, 0.5, 0.95]}
    rb_nb = show_performance_with_gsvc('NaiveBayes', MultinomialNB(), tfidf_x_train, y_train, params)
    print(rb_nb)

```

```
3564201863, 'train_score': 0.9966413429851532, 'val_score': 0.9056920458347125, 'best_model': MultinomialNB(alp
na=0.03))

show_metrics(rs_lr, rs_svc, rs_nb)
```

Model	train time (s)	val time (s)
rs_lr	111.330054755932	0.03
rs_svc	111.330054755932	0.03
rs_nb	111.330054755932	0.03

The top chart displays the scores for each model by train and validation sets. The y-axis represents the score, ranging from 0 to 40. The x-axis lists the models: Logistic Regression, SVC, and Naïve Bayes. The legend indicates that blue bars represent the train score and red bars represent the val score.

Model	Train Score	Val Score
Logistic Regression	38	0.04347893227471245
SVC	40	34.86272047625648
Naïve Bayes	0.05462815364201863	0.2047151062558492

The bottom chart displays the scores for each model by train and validation sets. The y-axis represents the score, ranging from 0.6 to 1.0. The x-axis lists the models: Logistic Regression, SVC, and Naïve Bayes. The legend indicates that blue bars represent the train score and red bars represent the val score.

Model	Train Score	Val Score
Logistic Regression	0.9519719569556	0.9040129170025326
SVC	0.9992487273953168	0.9067513114195242
Naïve Bayes	0.996641342881532	0.9125

Model	Dataset 1 (Blue)	Dataset 2 (Red)	Dataset 3 (White)
LogisticRegression	~0.45	~0.45	~0.45
SVC	~0.45	~0.45	~0.45
NaiveBayes	~0.45	~0.45	~0.45


```
[149] words_count = []
      for i in x_train:
          words_count.append(len(text_to_word_sequence(i, split=" ")))
      print(np.max(words_count))
      print(np.min(words_count))
      # 18- 26333
      NUM_WORDS = 20000
      MAX_LEN= 1000
      EMBEDDING_DIM = 100

      tokenizer = Tokenizer(num_words=NUM_WORDS)
      tokenizer.fit_on_texts(x_sp_train)
      word_index = tokenizer.word_index
      x_sp_train_dl = pad_sequences(tokenizer.texts_to_sequences(x_sp_train), maxlen=MAX_LEN)
      y_sp_train_dl = to_categorical(np.asarray(y_sp_train), num_classes=20)

      x_sp_val_dl = pad_sequences(tokenizer.texts_to_sequences(x_sp_val), maxlen=MAX_LEN)
      y_sp_val_dl = to_categorical(np.asarray(y_sp_val), num_classes=20)

      x_test_dl = pad_sequences(tokenizer.texts_to_sequences(x_test), maxlen=MAX_LEN)
      y_test_dl = to_categorical(np.asarray(y_test), num_classes=20)

16161
16

In [169]: # Finally using word2vec
          min_num_words = min(NUM_WORDS, len(word_index))
          embedding_matrix = np.zeros((min_num_words+1, EMBEDDING_DIM))

          for word, index in word_index.items():
              for word in words:
                  if word in model.wv:
                      vec = model.wv[word]
                      embedding_matrix[index] = model[word]

          print('embedding matrix shape: {}'.format(embedding_matrix.shape))
          embedding_layer = Embedding(min_num_words+1, EMBEDDING_DIM, weights=[embedding_matrix], input_length=MAX_LEN,
                                     trainable=False)
          text_cnn = text_CNN(embedding_layer)
          text_cnn.summary()

          embedding_matrix shape: (20001, 100)
          Model: "model_1"
```

```
In [169]: early_stopping = EarlyStopping(monitor='val_acc', patience=2, mode='max')

          history = text_cnn.fit(x_sp_train_dl, y_sp_train_dl, validation_data=(x_sp_val_dl, y_sp_val_dl), epochs=60, bat

Epoch 1/60
182/182 [=====] - 32s 173ms/step - loss: 2.9957 - acc: 0.0547 - val_loss: 2.9957 - val
acc: 0.0437
Epoch 2/60
182/182 [=====] - 33s 183ms/step - loss: 2.9957 - acc: 0.0547 - val_loss: 2.9957 - val
acc: 0.0437
Epoch 3/60
182/182 [=====] - 33s 184ms/step - loss: 2.9957 - acc: 0.0547 - val_loss: 2.9957 - val
acc: 0.0437

In [171]: show_history(history)
```



As you can see there cannot be done a qualitative work in such short time, and the results are corresponding to that fact.

I stopped further analysis as I realized that even to finish word2vec using code & explanation from kaggle() I would need a couple days.

And I assume this is fair for every other approach. Asking to do this in a week from people that most probably work 8 hours a day is cruel to say the least.

Sorry, I don't have time and resources to properly work with all of this. As a conclusion this happened to be the work that I physically and mentally cannot do.